

Performance measures for the two-node queue with finite buffers

Yanting Chen^a, Richard J. Boucherie^a, Jasper Goseling^a

^a*Stochastic Operations Research, University of Twente, P.O. Box 217, 7500 AE Enschede, The Netherlands*

Abstract

We consider two-node queue with finite buffers modeled as a two-dimensional random walk on a finite state space. We develop an approximation scheme based on the Markov reward approach to error bounds in order to bound performance measures of a two-dimensional finite random walk in terms of a perturbed random walk in which only the transitions along the boundaries are different from those in the original model. The invariant measure of the perturbed random walk is of product-form. We first apply this approximation scheme to a tandem queue with finite buffers and some variants of this model. Then, we show that our approximation scheme is sufficiently general by applying it to a coupled-queue with finite buffers and processor sharing.

Keywords: Random walk, Finite state space, Product-form, Error bounds, Performance measure

The two-node queue with a finite buffer space is one of the most extensively studied topics in queueing theory. Finite capacity buffers are usually used to model stochastic systems with limited storage capacity such as manufacturing, telecommunications or transportation applications. The two-node queue with a finite buffer space is often modeled as a two-dimensional random walk on a finite state space. Hence, it is sufficient to find performance measures of the corresponding two-dimensional random walk on a finite state space if we are interested in the performance of the two-node queue with a finite buffer space.

Email addresses: `y.chen@utwente.nl` (Yanting Chen), `r.j.boucherie@utwente.nl` (Richard J. Boucherie), `j.goseling@utwente.nl` (Jasper Goseling)

A special case of the two-node queue with a finite buffers space, which has been extensively studied so far, is the tandem queue with a finite buffer space. An extensive list of papers on this topic is provided in [2, 11]. Most of these papers focus on the development of approximations or algorithmic procedures to find steady-state system performances such as throughput and the average number of customers in the system. A popular approach used in such approximations is decomposition, see [1, 6]. The main variations of such a model are: three or more stations in the tandem queue [12], multiple servers at each station [20, 22], optimal design for allocating finite buffers to the stations [9], general service time [13, 18], etc. Numerical results of such approximations often imply that the proposed approximations are indeed the bounds of the specific performance measure. However, these approximation methods cannot be easily extended to a general method, which determines the steady-state performance measures of any two-node queue with finite buffers.

Van Dijk et al. [17] pioneered in developing error bounds for the system throughput using the product-form modification approach. The method has since been further developed by van Dijk et al. [15, 19] and has been applied to, for instance, Erlang loss networks [3], to networks with breakdowns [14], to queueing networks with non-exponential service [18] and to wireless communication networks with network coding [7]. An extensive description and overview of various applications of this method can be found in [16].

A major disadvantage of the error bound method mentioned above is that the verification steps that are required to apply the method can be technically quite complicated. Goseling et al. [8] developed a general verification technique for random walks in the quarter-plane. This verification technique is based on formulating the application of the error bounds method as solving a linear program. In doing so, it avoids completely the induction proof required in [19]. Moreover, instead of only bounding performance measures for specific queueing system, the approximation method developed in [8] accepts any random walk in the quarter-plane as an input.

The main contribution of the current work is to provide an approximation scheme which can be readily applied to approximate performance measures for any two-node queue with finite buffers. This is based on modifying the general verification technique developed in [8] for a two-dimensional random walk on a finite state space.

We apply this approximation scheme to a tandem queue with finite buffers. We show that the error bounds for the blocking probability are improved com-

pared with the error bounds for the blocking probability provided in [17]. The method in [17] is based on specific model modification. Apart from this, our approximation scheme is more general in the sense that other interesting performance measures could also be obtained easily. This is an advantage over the methods used in [15, 17, 19] where different model modifications are necessary for different performance measures. Moreover, we show that the error bounds can still be achieved for variations of the tandem queue with finite buffers. In particular, we consider the case that one server speeds-up or slows-down when another server is idle or saturated. Finally, we show that this approximation scheme also works for other two-node queue with finite buffer space model, for instance, a coupled-queue with processor sharing and finite buffers. The numerical results illustrate that our approximation scheme achieves tight bounds. Hence, our method can be useful for quick engineering purposes and optimal design of the queueing system.

The remainder of this paper proceeds as follows. In Section 1, we present the model and formulate the research problem. In Section 2, we provide an approximation scheme to bound performance measures for any two-node queue with finite buffers. We bound performance measures for a tandem queue with finite buffers and some variants of this model in Section 3. In Section 4, the approximation scheme has also been applied to a coupled-queue with processor sharing and finite buffers. Finally, we provide concluding remarks in Section 5.

1. Model and problem formulation

1.1. Two-node queue with finite buffers

The two-node queue with finite buffers is a queueing system with two servers, each of them having finite storage capacity. If a job arrives at a server which does not have any more storage capacity, then the job is lost. In general, the two queues influence each other, *i.e.*, the service rate at one of the queues depends on the number of jobs at the other.

Such a queueing system is naturally modeled as a two-dimensional finite random walk, which we introduce next. The connection between the continuous-time queueing system and the discrete-time random walk, obtained through uniformization, is made explicit for various examples in Section 3 and Section 4.

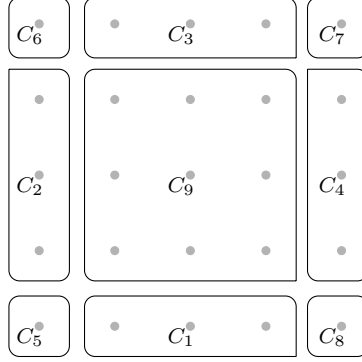


Figure 1: C -partition of S with components C_1, C_2, \dots, C_9 .

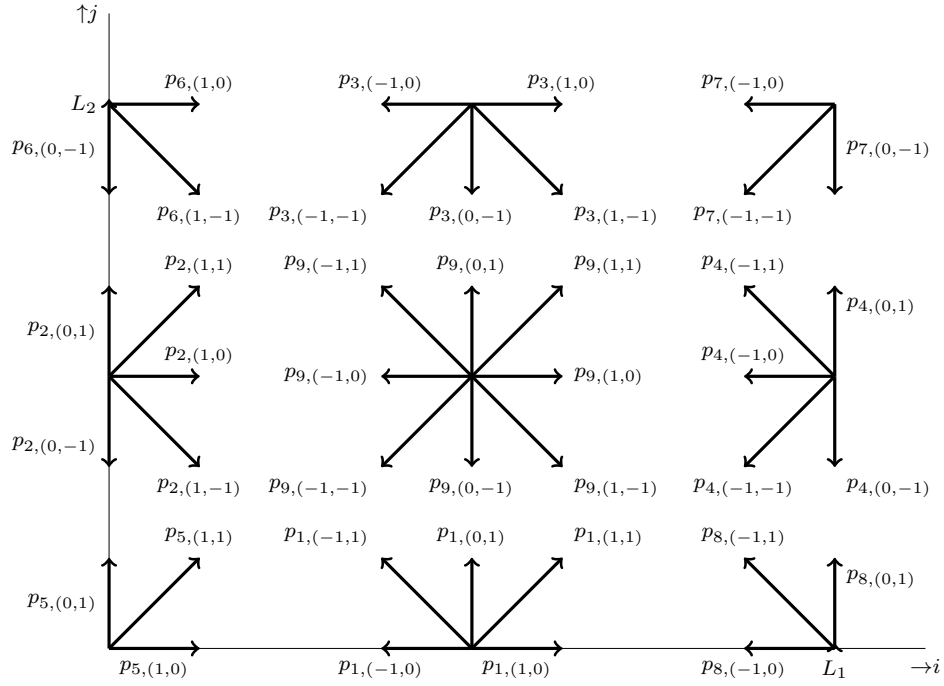


Figure 2: Two-dimensional finite random walk on S .

1.2. Two-dimensional finite random walk

We consider a two-dimensional random walk R on S where

$$S = \{0, 1, 2, \dots, L_1\} \times \{0, 1, 2, \dots, L_2\}.$$

We use a pair of coordinates to represent a state, *i.e.*, for $n \in S$, $n = (i, j)$. The state space is naturally partitioned in the following components (see Figure 1):

$$\begin{aligned} C_1 &= \{1, 2, 3, \dots, L_1 - 1\} \times \{0\}, & C_2 &= \{0\} \times \{1, 2, 3, \dots, L_2 - 1\}, \\ C_3 &= \{1, 2, 3, \dots, L_1 - 1\} \times \{L_2\}, & C_4 &= \{L_1\} \times \{1, 2, 3, \dots, L_2 - 1\}, \\ C_5 &= \{(0, 0)\}, & C_6 &= \{(0, L_2)\}, & C_7 &= \{(L_1, L_2)\}, & C_8 &= \{(L_1, 0)\}, \\ C_9 &= \{1, 2, 3, \dots, L_1 - 1\} \times \{1, 2, 3, \dots, L_2 - 1\}. \end{aligned}$$

We refer to this partition as the C -partition. The index of the component of state $n \in S$ is denoted by $k(n)$, *i.e.*, $n \in C_{k(n)}$. Transitions are restricted to the neighboring points (horizontally, vertically and diagonally). For $k = 1, 2, \dots, 9$, we denote by N_k the neighbors of a state in C_k . More precisely, $N_1 = \{-1, 0, 1\} \times \{0, 1\}$, $N_2 = \{0, 1\} \times \{-1, 0, 1\}$, $N_3 = \{-1, 0, 1\} \times \{-1, 0\}$, $N_4 = \{-1, 0\} \times \{-1, 0, 1\}$, $N_5 = \{0, 1\} \times \{0, 1\}$, $N_6 = \{0, 1\} \times \{-1, 0\}$, $N_7 = \{-1, 0\} \times \{-1, 0\}$, $N_8 = \{-1, 0\} \times \{1, 0\}$ and $N_9 = \{-1, 0, 1\} \times \{-1, 0, 1\}$. Also, let $N = N_9$.

Let $p_{k,u}$ denote the transition probability from state n in component k to $n + u$, where $u \in N_k$. The transition diagram of a two-dimensional finite random walk can be found in Figure 2. The system is homogeneous in the sense that the transition probabilities (incoming and outgoing) are translation invariant in each of the components, *i.e.*,

$$p_{k(n-u),u} = p_{k(n),u}, \quad \text{for } n - u \in S \text{ and } u \in k(n). \quad (1)$$

Equation (1) not only implies that the transition probabilities for each part of the state space are translation invariant but also ensures that also the transition probabilities entering the same component of the state space are translation invariant.

We assume that the random walk R that we consider is aperiodic, irreducible, positive recurrent, and has invariant probability measure $m(n)$, where $m(n)$ satisfies for all $n \in S$,

$$m(n) = \sum_{u \in N_{k(n)}} p_{k(n+u),-u} m(n+u).$$

1.3. Problem formulation

Our goal is to approximate the steady-state performance of the random walk R . The performance measure of interest is

$$\mathcal{F} = \sum_{n \in S} m(n) F(n),$$

where $F(n) : S \rightarrow [0, \infty)$ is linear in each of the components from C -partition, *i.e.*,

$$F(n) = f_{k(n),0} + f_{k(n),1}i + f_{k(n),2}j, \quad \text{for } n = (i, j) \in S. \quad (2)$$

The constants $f_{k(n),0}$, $f_{k(n),1}$ and $f_{k(n),2}$ are allowed to be different for different components from the C -partition of S .

In general, it is not possible to obtain the probability measure $m(n)$ in a closed-form. Therefore, we will use a perturbed random walk of which the invariant measure has a closed-form expression to approximate the performance measure \mathcal{F} .

We approximate the performance measure \mathcal{F} in terms of the perturbed random walk \bar{R} . We consider the perturbed random walk \bar{R} in which only the transition probabilities along the boundaries (C_1, \dots, C_8) are allowed to be different, *i.e.*, for instance, $p_{1,(-1,0)}$, $p_{1,(1,0)}$, $p_{1,(0,0)}$ for the state from C_1 are allowed to be different in \bar{R} , $p_{2,(0,1)}$, $p_{2,(0,-1)}$, $p_{2,(0,0)}$ for the state from C_2 are allowed to be different in \bar{R} , etc. An example of a perturbed random walk \bar{R} can be found in Figure 3.

We use $\bar{p}_{k,u}$ to denote the probability of \bar{R} jumping from any state n in component C_k to $n + u$, where $u \in N_k$. Moreover, let $q_{k,u} = \bar{p}_{k,u} - p_{k,u}$. The probability measure \bar{m} of \bar{R} is assumed to be of product-form, *i.e.*,

$$\bar{m}(n) = \alpha \rho^i \sigma^j,$$

where $n = (i, j)$ for some $(\rho, \sigma) \in (0, 1)^2$ and $\alpha \neq 0$. The measure \bar{m} is the invariant measure of \bar{R} , *i.e.*, it satisfies

$$\bar{m}(n) = \sum_{u \in N_{k(n)}} p_{k(n+u),-u} \bar{m}(n+u), \quad (3)$$

for all $n \in S$.

In the following sections, we are going to find upper and lower bounds of \mathcal{F} in terms of the perturbed random walk \bar{R} defined above.

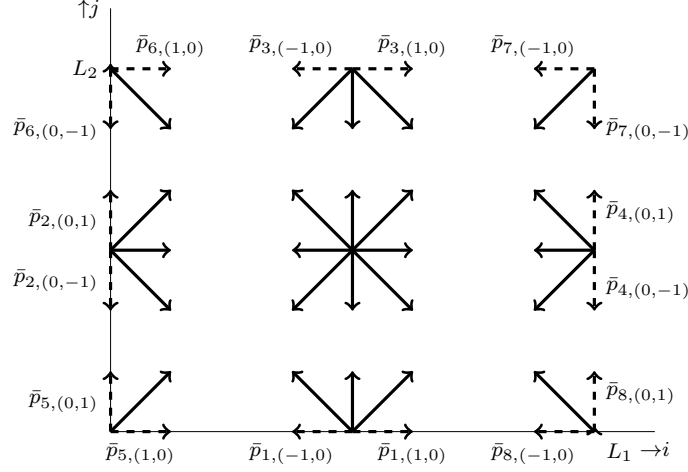


Figure 3: Perturbed random walk \bar{R} .

2. Proposed approximation scheme

In this section, we establish an approximation scheme to find upper and lower bounds for performance measures of a two-dimensional finite random walk.

In [8], an approximation scheme based on a linear program is developed for a random walk in the quarter-plane. This approximation scheme has also been used in [4]. We will show in this paper that the technique can be extended to cover our model, *i.e.*, a two-dimensional finite random walk. We will explain how this is achieved in the following sections.

2.1. Markov reward approach to error bounds

The fact that R and \bar{R} differ only along the boundaries of S makes it possible to obtain the error bounds for the performance measures via the Markov reward approach. An introduction to this technique is provided in [16]. We interpret F as a reward function, where $F(n)$ is the one step reward if the random walk is in state n . We denote by $F^t(n)$ the expected cumulative reward at time t if the random walk starts from state n at time 0, *i.e.*,

$$F^t(n) = \begin{cases} 0, & \text{if } t = 0, \\ F(n) + \sum_{u \in N_{k(n)}} p_{k(n),u} F^{t-1}(n+u), & \text{if } t > 0, \end{cases}$$

For convenience, let $F^t(n+u) = 0$ where $u \in \{(s,t) | s, t \in \{-1, 0, 1\}\}$ if $n+u \notin S$. Terms of the form $F^t(n+u) - F^t(n)$ play a crucial role in the Markov reward approach and are denoted as *bias terms*. Let $D_u^t = F^t(n+u) - F^t(n)$. For the unit vectors $e_1 = (1, 0)$, $e_2 = (0, 1)$, let $D_1^t(n) = D_{e_1}^t(n)$ and $D_2^t(n) = D_{e_2}^t(n)$.

The next result in [16] provides bounds for the approximation error for \mathcal{F} . We will use two non-negative functions \bar{F} and G to bound the performance measure \mathcal{F} .

Theorem 1 ([16]). *Let $\bar{F}: S \rightarrow [0, \infty)$ and $G: S \rightarrow [0, \infty)$ satisfy*

$$\left| \bar{F}(n) - F(n) + \sum_{u \in N_{k(n)}} q_{k(n),u} D_u^t(n) \right| \leq G(n), \quad (4)$$

for all $n \in S$ and $t \geq 0$. Then

$$\sum_{n \in S} [\bar{F}(n) - G(n)] \bar{m}(n) \leq \mathcal{F} \leq \sum_{n \in S} [\bar{F}(n) + G(n)] \bar{m}(n). \quad (5)$$

2.2. A linear program approach

In this section we present a linear program approach to bound the errors. Due to our construction of \bar{R} , the random walks R and \bar{R} differ only in the transitions that are along the unit directions, *i.e.*,

$$q_{k,u} = \bar{p}_{k,u} - p_{k,u} = 0 \quad \text{for } u \neq \{e_1, e_2, -e_1, -e_2, (0, 0)\}. \quad (6)$$

This restriction will significantly simplify the presentation of the result.

To start, consider the following optimization problem. We only consider how to obtain the upper bound for \mathcal{F} here because the lower bound for \mathcal{F} can be found similarly.

Problem 1

$$\text{minimize } \sum_{n \in S} [\bar{F}(n) + G(n)] \bar{m}(n), \quad (7)$$

$$\text{subject to } \left| \bar{F}(n) - F(n) + \sum_{s=1,2} (q_{k(n),e_s} D_s^t(n) + q_{k(n),-e_s} D_s^t(n - e_s)) \right|$$

$$\leq G(n), \quad \text{for } n \in S, t \geq 0, \quad (8)$$

$$\bar{F}(n) \geq 0, G(n) \geq 0, \quad \text{for } n \in S. \quad (9)$$

The variables in Problem 1 are the functions $\bar{F}(n)$, $G(n)$ and the parameters are $F(n)$, $\bar{m}(n)$, $q_{k(n),e_s}$ and $D_s^t(n)$ for $n \in S$, $s = 1, 2$. Hence, Problem 1 is a linear programming problem over two non-negative variables $\bar{F}(n)$ and $G(n)$ for every $n \in S$.

This linear program has infinitely many constraints because we have unbounded time horizon. We will first bound the bias term $D_s^t(n)$ uniformly over t . Then we have a linear program with a finite number of variables and constraints. However, further reduction is still needed because the number of variables and constraints will increase rapidly if L_1 and L_2 , which define the size of the state space, increase. Our contribution is to reduce Problem 1 to a linear programming problem where the number of variables and constraints does not depend on the size of the finite state space.

We now verify that the objective in Problem 1 is indeed an upper bound on the performance measure \mathcal{F} . Consider $D_{(0,0)}^t(n) = 0$, $D_{-e_s}^t(n) = -D_{e_s}^t(n - e_s)$ for $s = 1, 2$ and (6), it follows directly that (8) is equivalent to (4). Therefore, it follows from Theorem 1 that the objective of Problem 1 provides an upper bound on \mathcal{F} .

2.3. Bounding the bias terms

The main difficulty in solving Problem 1 is the unknown bias terms $D_s^t(n)$. It is in general not possible to find closed-form expressions for the bias terms. Therefore, we introduce two functions $A_s: S \rightarrow [0, \infty)$ and $B_s: S \rightarrow [0, \infty)$, $s = 1, 2$. We will formulate a finite number of constraints on functions A_s and B_s where $s = 1, 2$ such that for any t and $s = 1, 2$ we have

$$-A_s(n) \leq D_s^t(n) \leq B_s(n), \quad (10)$$

i.e., the functions A_s and B_s provide bounds on the bias terms uniformly over all $t \geq 0$. In the next section, we will find a finite number of constraints that imply (10). Our method is based on the method that was developed in [8] for the case of an unbounded state space.

For notational convenience, as will become clear below, we define a finer partition of S , the Z -partition. This partition is depicted in Figure 4. For example, we have $Z_1 = \{(0, 0)\}$, $Z_2 = \{(1, 0)\}$, $Z_3 = \{2, \dots, L_1 - 2\} \times \{0\}$, $Z_4 = \{(L_1 - 1, 0)\}$ and $Z_5 = \{(L_1, 0)\}$, the rest of the elements in the partition are determined similarly. Let $k^z(n)$ denote the label of component from Z -partition of state $n \in S$, *i.e.*, $n \in Z_{k^z(n)}$. Similar to the definition of N_k , let N_k^z denote the neighbors of a state n in Z_k from the Z -partition of S .

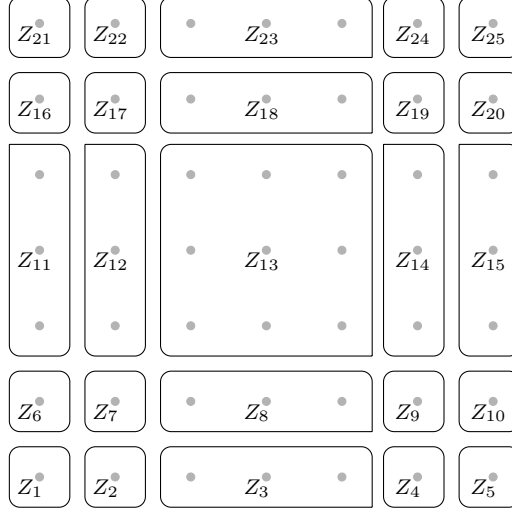


Figure 4: Z -partition of S with components Z_1, Z_2, \dots, Z_{25} .

The constraints which ensure (10) are obtained based on an induction in t . More precisely, we express D_s^{t+1} as a linear combination of D_1^t and D_2^t as

$$D_s^{t+1}(n) = F(n + e_s) - F(n) + \sum_{v=1,2} \sum_{u \in N_{k(n)}^z} c_{s,k^z(n),v,u} D_v^t(n + u), \quad (11)$$

where the $c_{s,k,v,u}$, $s \in \{1, 2\}$, $k \in \{1, 2, \dots, 25\}$, $v \in \{1, 2\}$, $u \in N_k^z$ are constants. An important property of the Z -partition is that starting from any state n in component k^z of the Z -partition the component $k(n + u)$ in the C -partition is well defined for all $u \in N_k^z$ and depends only on k^z and u . In [8] it was shown, using this property, that constants $c_{s,k,v,u}$ that ensure (11) always exist and that they can be expressed as simple functions of the transition probabilities of the random walk. The results in [8] are derived for the random walk on the whole quarter-plane. However, a careful inspection of the results in [8] reveals that they hold also for our model of a random walk on a bounded state space. Therefore, we refer the reader to [8] and omit further details here.

We are now ready to bound the bias terms based on (11). The result, which is easy to verify, states that if $A_s: S \rightarrow [0, \infty)$ and $B_s: S \rightarrow [0, \infty)$

where $s = 1, 2$ satisfy

$$F(n + e_s) - F(n) + \sum_{v=1,2} \sum_{u \in N_{k(n)}^z} \max\{-c_{s,k^z(n),v,u} A_s(n+u), c_{s,k^z(n),v,u} B_s(n+u)\} \leq B_s(n),$$

$$F(n) - F(n + e_s) + \sum_{v=1,2} \sum_{u \in N_{k(n)}^z} \max\{-c_{s,k^z(n),v,u} B_s(n+u), c_{s,k^z(n),v,u} A_s(n+u)\} \leq A_s(n),$$

for all $n \in S$, then

$$-A_s(n) \leq D_s^t(n) \leq B_s(n),$$

for $s = 1, 2$, $n \in S$ and $t \geq 0$.

After bounding the bias terms, we are able to rewrite the linear program Problem 1 into Problem 2 with a new variables $E_s(n)$ where $s = 1, 2$ and $n \in S$.

Problem 2

$$\text{minimize} \quad \sum_{n \in S} [\bar{F}(n) + G(n)] \bar{m}(n),$$

$$\begin{aligned} \text{subject to} \quad & \left| \bar{F}(n) - F(n) + \sum_{s=1,2} (q_{k(n),e_s} E_s(n) + q_{k(n),-e_s} E_s(n - e_s)) \right| \\ & \leq G(n), \\ & -A_s(n) \leq E_s(n) \leq B_s(n), \\ & F(n + e_s) - F(n) \\ & + \sum_{v=1,2} \sum_{u \in N_{k(n)}^z} \max\{-c_{s,k^z(n),v,u} A_s(n+u), c_{s,k^z(n),v,u} B_s(n+u)\} \\ & \leq B_s(n), \\ & F(n) - F(n + e_s) \\ & + \sum_{v=1,2} \sum_{u \in N_{k(n)}^z} \max\{-c_{s,k^z(n),v,u} B_s(n+u), c_{s,k^z(n),v,u} A_s(n+u)\} \\ & \leq A_s(n), \\ & \bar{F}(n) \geq 0, G(n) \geq 0, A_s(n) \geq 0, B_s(n) \geq 0, \\ & \text{for } n \in S, s \in \{1, 2\}. \end{aligned}$$

2.4. Fixed number of variables and constraints

The final step is to reduce Problem 2 to a linear program with fixed number of variables and constraints regardless of the size of the state space.

We first introduce the notion of a piecewise-linear function on the Z -partition. A function $F : S \rightarrow [0, \infty)$ is called Z -linear if the function is linear in each of the components from Z -partition, *i.e.*,

$$F(n) = f_{k^z(n),0} + f_{k^z(n),1}i + f_{k^z(n),2}j, \quad \text{for } n = (i, j) \in S.$$

where $f_{k^z(n),0}$, $f_{k^z(n),1}$ and $f_{k^z(n),2}$ are the constants that define the function. In similar fashion we define C -linear functions on the C -partition of S .

Now, in Problem 2 we put the additional constraint that the variables \bar{F} , G , A_s , B_s and E_s are C -linear functions. Hence, these functions are defined in terms of variables, the number of which is independent on L_1 and L_2 . Hence, the number of variables in the resulting linear program is independent of L_1 and L_2 .

It remains to show that the number of constraints is independent of L_1 and L_2 . Following the reasoning on the properties of Z -partition below (11) it is easy to see that all constraints in Problem 2 can be formulated as a non-negativity constraint on a Z -linear function. Such a constraint on a Z -linear function induces at most 4 constraints per component in the Z -partition, one constraint for each corner of the component. This indicates that the number of constraints does not depend on the size of the state space, since the number of constraints are fixed as well.

2.5. The optimal solutions

We are now able to find the upper and lower bounds of \mathcal{F} based on the linear program here.

Let \mathcal{P} denote the set of (\bar{F}, G) for which we are able to find functions A_s , B_s and E_s where $s = 1, 2$ such that all constraints in Problem 2 are satisfied. Then, we find the upper and lower bounds for \mathcal{F} as follow.

$$\mathcal{F}_{up} = \max \left\{ \sum_{n \in S} [\bar{F}(n) + G(n)] \bar{m}(n) \mid (\bar{F}, G) \in \mathcal{P} \right\},$$

and

$$\mathcal{F}_{low} = \min \left\{ \sum_{n \in S} [\bar{F}(n) - G(n)] \bar{m}(n) \mid (\bar{F}, G) \in \mathcal{P} \right\}.$$

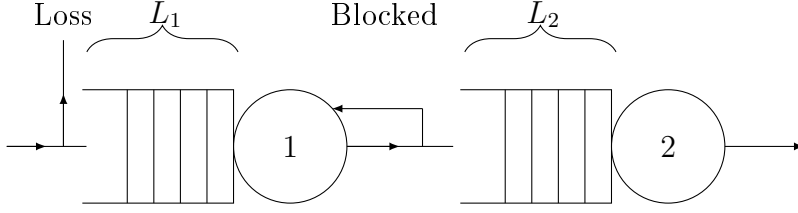


Figure 5: Tandem queue with finite buffers.

We have now presented the complete approximation scheme to obtain the upper and lower bounds for \mathcal{F} using the perturbed random walk \bar{R} of which the probability measure is of product-form.

In the following sections, we will consider two examples: a tandem queue with finite buffers and a coupled-queue with processor sharing and finite buffers.

3. Application to the Tandem queue with finite buffers

In this section, we investigate the applications of the approximation scheme proposed in Section 2.

3.1. Model description

Consider a two-node tandem queue with Poisson arrivals at rate λ . Both nodes have a single server. At most a finite number of jobs, say L_1 and L_2 jobs, can be present at nodes 1 and 2. This includes the jobs in service. An arriving job is rejected if node 1 is saturated, *i.e.*, there are L_1 jobs at node 1. The service time for the jobs at both nodes is exponentially distributed with parameters μ_1 and μ_2 , respectively.

When node 2 is saturated, *i.e.*, there are L_2 jobs at node 2, node 1 stops serving. When it is not blocked, it instantly routes to node 2. All service times are independent. We also assume that the service discipline is first-in first-out.

The tandem queue with finite buffers can be represented by a continuous-time Markov process whose state space consists of the pairs (i, j) where i and j are the number of jobs at node 1 and node 2, respectively. We now uniformize this continuous-time Markov process to obtain a discrete-time random walk. We assume without loss of generality that $\lambda + \mu_1 + \mu_2 \leq$

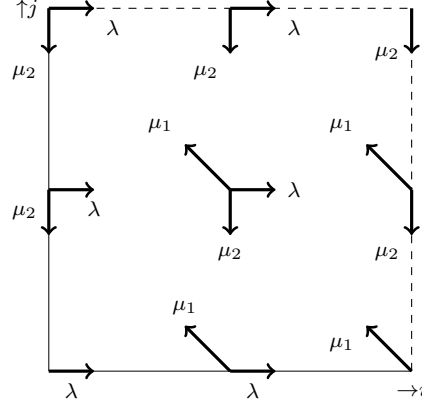


Figure 6: Transition diagram of R_T .

1 and uniformize the continuous-time Markov process with uniformization parameter 1. We denote this random walk by R_T . All transition probabilities of R_T , except those for the transitions from a state to itself, are illustrated in Figure 6.

3.2. Perturbed random walk of R_T

We now present a perturbed random walk \bar{R}_T . The invariant measure of the perturbed random walk \bar{R}_T is of product-form and only the transitions along the boundaries in \bar{R}_T are different from those in R_T .

In the perturbed random walk \bar{R}_T , the transition probabilities in the components C_3, C_4, C_6, C_7, C_8 are different from those in R_T . More precisely, we have $\bar{p}_{3,(1,0)} = \lambda$, $\bar{p}_{3,(-1,0)} = \mu_1$, $\bar{p}_{4,(0,1)} = \lambda$, $\bar{p}_{4,(0,-1)} = \mu_2$, see Figure 7. It can be readily verified that the measure, which is of product-form, with α , which depends on L_1 and L_2 as the normalizing constant

$$\bar{m}(i, j) = \alpha \left(\frac{\lambda}{\mu_1} \right)^i \left(\frac{\lambda}{\mu_2} \right)^j$$

is the probability measure of the perturbed random walk by substitution into the global balance equations (3) together with the normalization requirement.

3.3. Bounding the blocking probability

In this section, we provide error bounds for the blocking probability for the tandem queue with finite buffers using our approximation scheme provided in Section 2. Moreover, we show that our results are better than those obtain by van Dijk et al. in [17].

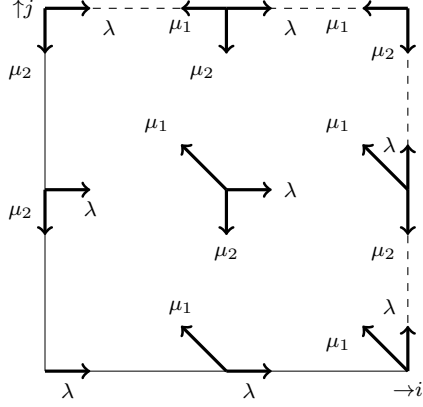


Figure 7: Transition diagram of \bar{R}_T .

For a given performance measure \mathcal{F} , we use \mathcal{F}^{up} , \mathcal{F}^{low} to denote the upper and lower bounds for \mathcal{F} obtained based on our approximation scheme and $\tilde{\mathcal{F}}^{up}$, $\tilde{\mathcal{F}}^{low}$ to denote the upper and lower bounds based on the method suggested by van Dijk et al. [17].

We use \mathcal{F}_0 to denote the blocking probability, *i.e.*, the probability that an arriving job is rejected. We now consider an example that has also been considered in [17].

Example 1. Consider a tandem queue with finite buffers, we have $\lambda = 0.1$, $\mu_1 = 0.2$, $\mu_2 = 0.2$.

We would like to compute the blocking probability of the queueing system. Hence, for the performance measure function $F(n)$, defined in (2), we set the coefficients $f_{k,d}$ where with $k = 1, 2, \dots, 9$, $d = 1, 2, 3$ to be $f_{8,1} = 1$, $f_{4,1} = 1$, $f_{7,1} = 1$ and others 0. The error bounds can be found in Figure 8. Clearly, our results outperform the error bounds obtained in [17]. Moreover, the difference between the upper and lower bounds of \mathcal{F}_0 are captured in Figure 9. This indicates that our error bounds are tighter than those in [17].

In addition to the improved bounds, there is another advantage to our method. There is a limitation to the model modification approach that is used in [17]. This method requires a different model modification for each specific performance measure. For instance, the specific model modifications which are used to find error bounds for the blocking probability of a tandem

Example 1
$\lambda = 0.1$
$\mu_1 = 0.2$
$\mu_2 = 0.2$
$L_1 = L_2$

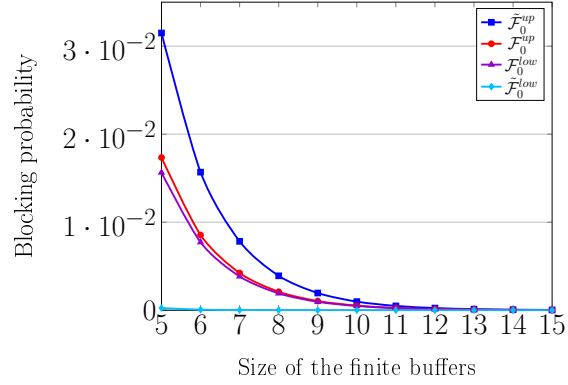


Figure 8: The blocking probability \mathcal{F}_0 .

Example 1
$\lambda = 0.1$
$\mu_1 = 0.2$
$\mu_2 = 0.2$
$L_1 = L_2$

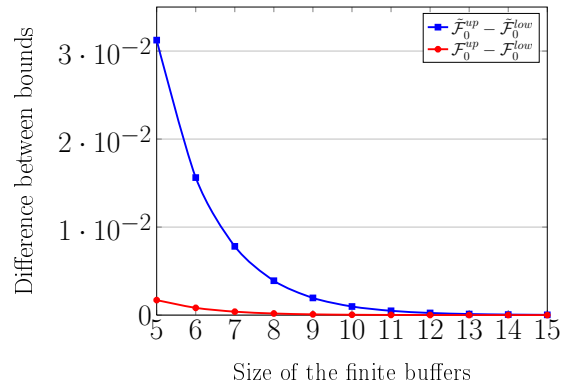


Figure 9: The difference between bounds of \mathcal{F}_0 .

queue with finite buffers in [17] cannot be used to obtain error bounds for the average number of jobs in the first node. In addition, extra effort is needed to verify that the model modifications are indeed valid for a specific performance measure. In the next section, we will show that our method can easily provide error bounds for other performance measures without extra effort.

3.4. Bounds for other performance measures

In this section, we will demonstrate the error bounds for other performance measures for Example 1, *i.e.*, a tandem queue with finite buffers.

Let \mathcal{F}_1 be the average number of jobs at node 1 and \mathcal{F}_2 which is the average number of jobs at node 2.

In general, the models, (*i.e.*, the perturbed systems), used to bound the blocking probability in [17] cannot be used to bound \mathcal{F}_1 and \mathcal{F}_2 . The method in [17] requires different upper and lower bound models for different performance measures. Moreover, this method also requires effort to verify that they are indeed the upper and lower bound models for this specific performance measure. Our approximation scheme does not have this disadvantage. For different performance measure, we only need to change the coefficients $f_{k,d}$ where $k = 1, 2, \dots, 9$ and $d = 1, 2, 3$ in $F(n)$, which is defined in (2).

It can be readily verified that the performance measure \mathcal{F} is \mathcal{F}_1 if and only if we assign following values to the coefficients: $f_{1,2} = 1, f_{8,2} = 1, f_{9,2} = 1, f_{4,2} = 1, f_{3,2} = 1, f_{7,2} = 1$ and others 0. Figure 10 presents the error bounds of \mathcal{F}_1 . Similarly, the performance measure \mathcal{F} is \mathcal{F}_2 if and only if we assign following values to the coefficients: $f_{2,3} = 1, f_{9,3} = 1, f_{4,3} = 1, f_{6,3} = 1, f_{3,3} = 1, f_{7,3} = 1$ and others 0. Figure 11 presents the error bounds of \mathcal{F}_2 .

The results show that tight bounds have been achieved with our approximation scheme. Moreover, the only thing we need to change for different performance measures is the input function, which does not require further model modifications. In the next section, we will show that our approximation scheme could also give error bounds for the performance measures of the tandem queue with finite buffers which has a slower or faster server when another node is idle or saturated, respectively, without model modifications as well.

3.5. Tandem queue with finite buffers and server slow-down/speed-up

In this section, we consider two variants of the tandem queue with finite buffers. More specifically, we provide error bounds for the blocking proba-

Example 1
$\lambda = 0.1$
$\mu_1 = 0.2$
$\mu_2 = 0.2$
$L_1 = L_2$

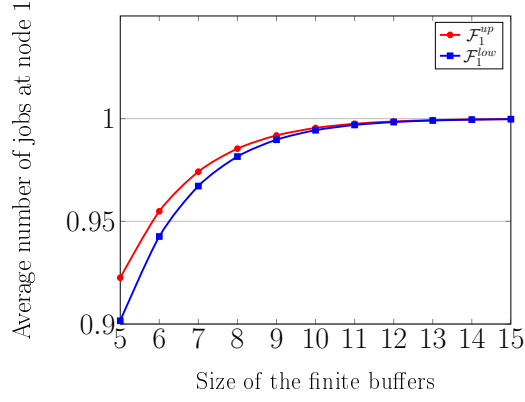


Figure 10: Average number of jobs at node 1, \mathcal{F}_1 .

Example 1
$\lambda = 0.1$
$\mu_1 = 0.2$
$\mu_2 = 0.2$
$L_1 = L_2$

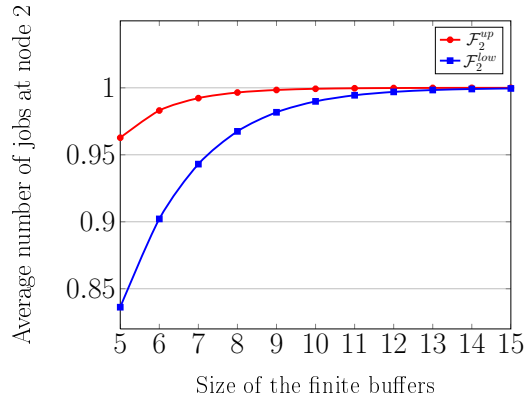


Figure 11: Average number of jobs at node 2, \mathcal{F}_2 .

bilities when one server in the tandem with finite buffers is slower or faster if another node is idle or saturated, respectively.

3.5.1. Tandem queue with finite buffers and server slow-down

Tandem queue with server slow-down has been previously studied in, for instance, [21, 10]. A specific type of tandem queue with finite buffers and server slow-down has been considered in [21, 10]. More precisely, the service speed of node 1 is reduced as soon as the number of jobs in node 2 reaches some pre-specified threshold because of some sort of protection against frequent overflows.

We consider a different scenario with server slow-down. In our case, the service rate at node 2 reduces when node 1 is idle. This comes from a practical situation that when node 1 is idle, the working pressure for node 2 decreases and can shift some working capacity to other tasks. Therefore, we consider a two-node tandem queue with Poisson arrivals at rate λ . Both nodes have a single server. At most a finite number of jobs, say L_1 and L_2 jobs, can be present at nodes 1 and 2, respectively. An arriving job is rejected if node 1 is saturated. The service time for the jobs at both nodes are exponentially distributed with parameters μ_1 and μ_2 , respectively. While node 2 is saturated, node 1 stops serving. When it is not blocked, it instantly routes to node 2. While node 1 is idle, the service rate of node 2 becomes $\tilde{\mu}_2$ where $\tilde{\mu}_2 < \mu_2$. All service times are independent. We also assume that the service discipline is first-in first-out.

The tandem queue with finite buffers and server slow-down can be represented by a continuous-time Markov process whose state space consists of the pairs (i, j) where i and j are the number of jobs at node 1 and node 2, respectively. We assume without loss of generality that $\lambda + \mu_1 + \mu_2 \leq 1$ and uniformize this continuous-time Markov process with uniformization parameter 1. Then we obtain a discrete-time random walk. We denote this random walk by R_T^{sd} , all transition probabilities of R_T^{sd} , except those for the transitions from a state to itself, are illustrated in Figure 12.

It can be readily verified that the random walk \bar{R}_T as defined in Section 3.2 is a perturbed random walk of R_T^{sd} as well, *i.e.*, the transition probabilities in \bar{R}_T only differ from those in R_T^{sd} along the boundaries. We next consider a numerical example.

Example 2 (slow-down). *Consider a tandem queue with finite buffers and server slow-down, we have $\lambda = 0.1$, $\mu_1 = 0.2$, $\mu_2 = 0.2$ and $\tilde{\mu}_2 = 0.5\mu_2$.*

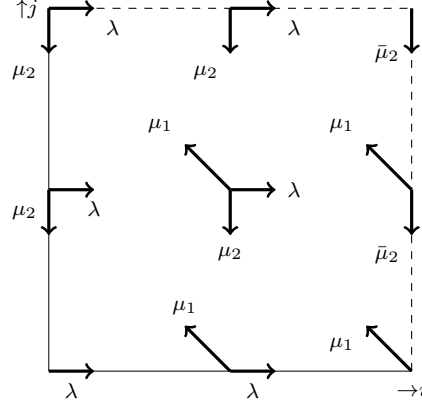


Figure 14: Tandem queue with finite buffers and server speed-up.

We consider the following scenario with server speed-up: The service rate at node 2 increases when node 1 is saturated. This comes from a practical situation, for instance, when node 1 is saturated, the working pressure for node 2 increases to eliminate the jobs in the queueing system. Therefore, we consider a two-node tandem queue with Poisson arrivals at rate λ . Both nodes have a single server. At most a finite number of jobs, say L_1 and L_2 jobs, can be present at nodes 1 and 2, respectively. An arriving job is rejected if node 1 is saturated. The service time for the jobs at both nodes are exponential distributed with parameters μ_1 and μ_2 , respectively. When node 2 is saturated, node 1 stops serving. When it is not blocked, it instantly routes to node 2. When node 1 is saturated, the service rate of node 2 becomes $\bar{\mu}_2$ where $\bar{\mu}_2 > \mu_2$. All service times are independent. We also assume that the service discipline is first-in first-out.

Tandem queue with finite buffers and server speed-up can be represented by a continuous-time Markov process whose state space consists of the pairs (i, j) where i and j are the number of jobs at node 1 and node 2, respectively. We assume without loss of generality that $\lambda + \mu_1 + \bar{\mu}_2 \leq 1$ and uniformize this continuous-time Markov process with uniformization parameter 1. Then we obtain a discrete-time random walk. We denote this random walk by R_T^{su} , all transition probabilities of R_T^{su} , except those for the transitions from a state to itself, are illustrated in Figure 14.

Again, it can be readily verified that the random walk \bar{R}_T as defined in Section 3.2 is a perturbed random walk of R_T^{su} because only the transitions along the boundaries in \bar{R}_T are different from those in R_T^{su} . We next consider

Example 3
$\lambda = 0.1$
$\mu_1 = 0.2$
$\mu_2 = 0.2$
$L_1 = L_2$
$\bar{\mu}_2 = 1.2\mu_2$

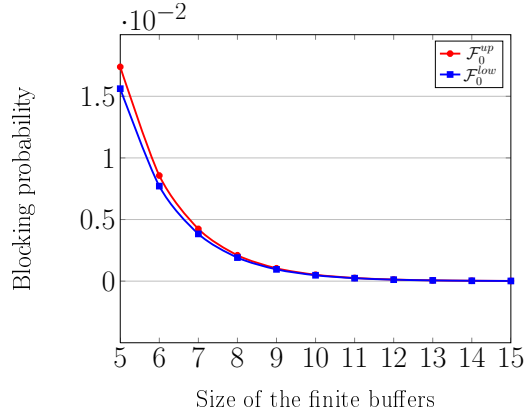


Figure 15: Blocking probability with server speed-up.

the following numerical example.

Example 3 (speed-up). Consider a tandem queue with finite buffers and server speed-up, we have $\lambda = 0.1$, $\mu_1 = 0.2$, $\mu_2 = 0.2$ and $\bar{\mu}_2 = 1.2\mu_2$.

The error bounds for the blocking probability of Example 3 can be found in Figure 15.

Until now, we have focused on the tandem queue with finite buffers model. In fact, our approximation scheme has been constructed in such a manner that it can be applied to any two-node queueing system with finite buffers. In the next section, we will study a different two-node queueing system with finite buffers and obtain error bounds for some performance measures.

4. Application to the coupled-queue with processor sharing and finite buffers

In this section, we apply the approximation scheme to a coupled-queue with processor sharing and finite buffers. Two coupled processors problem has been extensively studied so far. In particular, Fayolle et al. reduce the problem of finding the generating function of the invariant measure to a Riemann-Hilbert problem in [5]. However, when we have finite buffers, in general, the methods developed for a coupled-queue with infinite buffers are no longer valid.

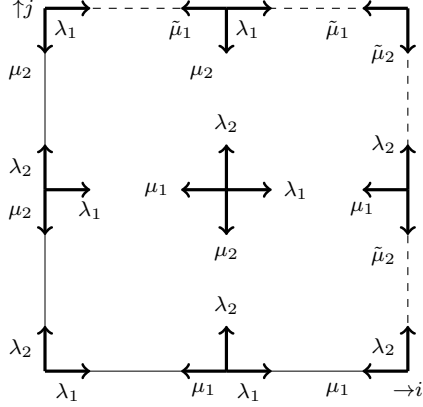


Figure 16: Coupled-queue with processor sharing and finite buffers.

4.1. Model description

Consider a two-node queue with Poisson arrivals at rate λ_1 for node 1 and λ_2 for node 2. Both nodes have a single server and at most L_1 and L_2 jobs can be present at nodes 1 and 2, respectively. When neither of the queues is saturated they evolve independently, but when one becomes saturated the service rate in the other changes. An arriving job for node 1 is rejected when node 1 is saturated. Similarly, an arriving job for node 2 is rejected when node 2 is saturated. The service time at both nodes is exponentially distributed with parameters μ_1 and μ_2 , respectively, when neither of the nodes is saturated. When node 1 is saturated, the service rate at node 2 becomes $\tilde{\mu}_2$ where $\tilde{\mu}_2 > \mu_2$ and when node 2 is saturated, the service rate at node 1 becomes $\tilde{\mu}_1$ where $\tilde{\mu}_1 > \mu_1$. All service requirements are independent. We also assume that the service discipline is first-in first-out.

This coupled-queue with processor sharing and finite buffers can be represented by a continuous-time Markov process whose state space consists of the pairs (i, j) where i and j are the number of jobs at node 1 and node 2, respectively. We assume without loss of generality that $\lambda_1 + \lambda_2 + \tilde{\mu}_1 + \tilde{\mu}_2 \leq 1$ and uniformize this continuous-time Markov process with uniformization parameter 1. Then we obtain a discrete-time random walk. We denote this random walk by R_C . All transition probabilities of R_C , except those for the transitions from a state to itself, are illustrated in Figure 16.

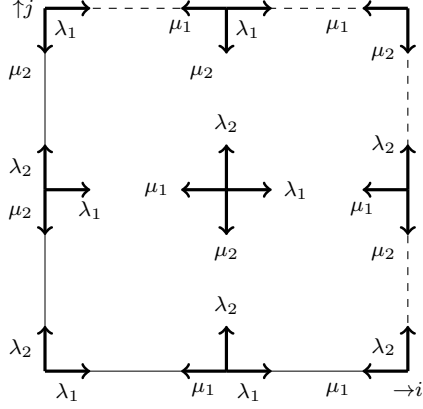


Figure 17: Transition diagram of the perturbed random walk \bar{R}_C .

4.2. Perturbed random walk \bar{R}_C

We now display a perturbed random walk \bar{R}_C of R_C such that the probability measure of \bar{R}_C is of product-form and only the transitions along the boundaries in \bar{R}_C are different from those in R_C .

In the perturbed random walk \bar{R}_C , the transition probabilities in components C_3, C_4, C_6, C_7, C_8 are different from that in R_C . More precisely, we have $p_{1,(1,0)} = \lambda_1, p_{1,(-1,0)} = \mu_1, p_{2,(0,1)} = \lambda_2, p_{2,(0,-1)} = \mu_2, p_{3,(1,0)} = \lambda_1, p_{3,(-1,0)} = \mu_1, p_{4,(0,1)} = \lambda_2, p_{4,(0,-1)} = \mu_2$, see Figure 17. It can be readily verified that the measure, which is of product-form, with α , which depends on L_1 and L_2 as the normalizing constant

$$\bar{m}(n) = \alpha \left(\frac{\lambda_1}{\mu_1} \right)^i \left(\frac{\lambda_2}{\mu_2} \right)^j \quad \text{where } n = (i, j),$$

is the probability measure of the perturbed random walk by substitution into the global balance equations (3) together with the normalization requirement.

We next illustrate a numerical example of a coupled-queue with processor sharing and finite buffers.

4.3. Numerical results

Example 4. Consider a coupled-queue with processor sharing and finite buffers, we have $\lambda_1 = \lambda_2 = 0.1, \mu_1 = \mu_2 = 0.2, \tilde{\mu}_1 = 0.4, \tilde{\mu}_2 = 0.3$.

Example 4
$\lambda_1 = \lambda_2 = 0.1$
$\mu_1 = \mu_2 = 0.2$
$L_1 = L_2$
$\tilde{\mu}_1 = 0.4$
$\tilde{\mu}_2 = 0.3$

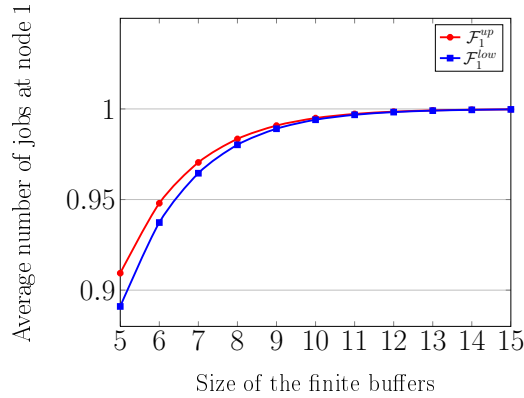


Figure 18: Average number of jobs at node 1, \mathcal{F}_1 .

For Example 4, the error bounds for the average number of jobs at node 1 can be found in Figure 18. This result also indicates that our approximation is quite general in the sense that it can be used to find the performance measure of another two-node queue with finite buffers, as, for example, a coupled-queue with processor sharing and finite buffers.

5. Conclusion

In this paper, we presented a generalized approximation scheme for a two-node queue with finite buffers that establishes error bounds for a large class of performance measures. Our work is an extension of the linear programming approach developed in [8] to approximate performance measures of random walks in the quarter-plane. We applied our approximation scheme to obtain bounds for performance measures of a tandem queue with finite buffers and some variants of this model. Then, we applied our approximation scheme to a coupled-queue with processor sharing and finite buffers. The approximation scheme gives tight bounds for various performance measures, like the blocking probability and the average number of jobs at node 1.

References

- [1] N. Asadathorn and X. Chao. A decomposition approximation for assembly-disassembly queueing networks with finite buffer and blocking. *Annals of Operations research*, 87:247–261, 1999.

- [2] S. Balsamo. Queueing networks with blocking: Analysis, solution algorithms and properties. In *Network performance engineering*, pages 233–257. Springer, 2011.
- [3] R.J. Boucherie and N.M. van Dijk. Monotonicity and error bounds for networks of erlang loss queues. *Queueing systems*, 62(1-2):159–193, 2009.
- [4] Y. Chen, R. J Boucherie, and J. Goseling. Invariant measures and error bounds for random walks in the quarter-plane based on sums of geometric terms. *arXiv preprint arXiv:1502.07218*, 2015.
- [5] G. Fayolle and R. Iasnogorodski. Two coupled processors: the reduction to a riemann-hilbert problem. *Probability Theory and Related Fields*, 47(3):325–351, 1979.
- [6] S. B. Gershwin. An efficient decomposition method for the approximate evaluation of tandem queues with finite storage space and blocking. *Operations Research*, 35(2):291–305, 1987.
- [7] J. Goseling, R.J. Boucherie, and J.C.W. van Ommeren. Energy–delay tradeoff in a two-way relay with network coding. *Performance Evaluation*, 70(11):981–994, 2013.
- [8] J. Goseling, R.J. Boucherie, and J.C.W van Ommeren. A linear programming approach to error bounds for random walks in the quarter-plane. *arXiv preprint arXiv:1409.3736*, 2014.
- [9] F.S. Hillier and K.C. So. On the optimal design of tandem queueing systems with finite buffers. *Queueing Systems*, 21(3-4):245–266, 1995.
- [10] D.I. Miretskiy, W.R.W. Scheinhardt, and M.R.H. Mandjes. State-dependent importance sampling for a slowdown tandem queue. *Annals of Operations Research*, 189(1):299–329, 2011.
- [11] H.G. Perros. *Queueing networks with blocking*. Oxford University Press, Inc., 1994.
- [12] J.G. Shanthikumar and M.A. Jafari. Bounding the performance of tandem queues with finite buffer spaces. *Annals of Operations Research*, 48(2):185–195, 1994.

- [13] N.M. van Dijk. A formal proof for the insensitivity of simple bounds for finite multi-server non-exponential tandem queues based on monotonicity results. *Stochastic processes and their applications*, 27:261–277, 1987.
- [14] N.M. van Dijk. Simple bounds for queueing systems with breakdowns. *Performance Evaluation*, 8(2):117–128, 1988.
- [15] N.M. van Dijk. Bounds and error bounds for queueing networks. *Annals of Operations Research*, 79:295–319, 1998.
- [16] N.M. van Dijk. Error bounds and comparison results: The Markov reward approach for queueing networks. In R.J. Boucherie and N.M. Van Dijk, editors, *Queueing Networks: A Fundamental Approach*, volume 154 of *International Series in Operations Research & Management Science*. Springer, 2011.
- [17] N.M. van Dijk and B.F. Lamond. Simple bounds for finite single-server exponential tandem queues. *Operations research*, pages 470–477, 1988.
- [18] N.M. van Dijk and M. Miyazawa. Error bounds for perturbing nonexponential queues. *Mathematics of Operations Research*, 29(3):525–558, 2004.
- [19] N.M. van Dijk and M.L. Puterman. Perturbation theory for Markov reward processes with applications to queueing systems. *Advances in Applied Probability*, 20(1):79–98, 1988.
- [20] N.M. van Dijk and J. van der Wal. Simple bounds and monotonicity results for finite multi-server exponential tandem queues. *Queueing Systems*, 4(1):1–15, 1989.
- [21] N.D. van Foreest, J.C.W. van Ommeren, M.R.H. Mandjes, and W.R.W. Scheinhardt. A tandem queue with server slow-down and blocking. *Stochastic Models*, 21(2-3):695–724, 2005.
- [22] M. van Vuuren, I.J.B.F. Adan, and S.A.E. Resing-Sassen. Performance analysis of multi-server tandem queues with finite buffers and blocking. In *Stochastic Modeling of Manufacturing Systems*, pages 169–192. Springer, 2006.